

## LOAD DISTRIBUTION CONTROL SYSTEM AND ITS DEVICE

**Publication number:** JP11312149 (A)

**Publication date:** 1999-11-09

**Inventor(s):** HIRATSUKA MASASHI

**Applicant(s):** HITACHI LTD

**Classification:**

- **international:** G06F15/16; G06F9/50; G06F15/177; G06F15/16; G06F9/46; (IPC1-7): G06F15/16

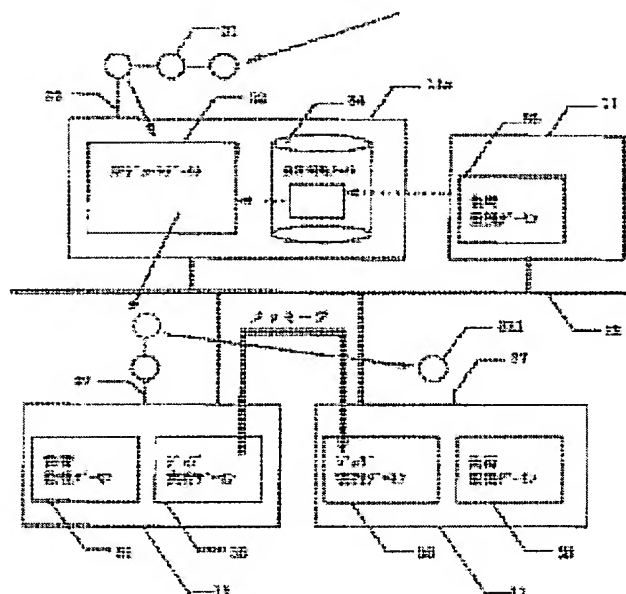
- **European:**

**Application number:** JP19980118122 19980428

**Priority number(s):** JP19980118122 19980428

### Abstract of JP 11312149 (A)

**PROBLEM TO BE SOLVED:** To keep the balance of loads among computers by adding load information which indicates the load state of the self computer to a transmission message at the time of an inter-node communication processings among the optional computers and exchanging the load information between the computers. **SOLUTION:** Load information is transmitted from a load monitor demon 35 on the computer 11 to the computer 11a by a certain opportunity and is stored in plural pieces of load information file 34. Then, a scheduler demon 33 on the computer 11a determines by which one of the computers 11 a job or a process 31, which reaches a reception queue 32 is executed based on the pieces of load information, queues the job or the process 31 in the execution waiting queue 37 of the applying computer 11, successively takes-out it by a job execution demon and exchanges the pieces of load information between the respective computers 11. Thus, the balance of the loads among the respective computers 11 is kept.



Data supplied from the **esp@cenet** database — Worldwide

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-312149

(43) 公開日 平成11年(1999)11月9日

(51) Int.Cl.<sup>6</sup>

G 0 6 F 15/16

識別記号

3 8 0

F I

G 0 6 F 15/16

3 8 0 Z

審査請求 未請求 請求項の数2 O L (全 5 頁)

(21) 出願番号 特願平10-118122

(22) 出願日 平成10年(1998)4月28日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 平塚 正史

神奈川県横浜市戸塚区戸塚町5030番地 株

式会社日立製作所ソフトウェア開発本部内

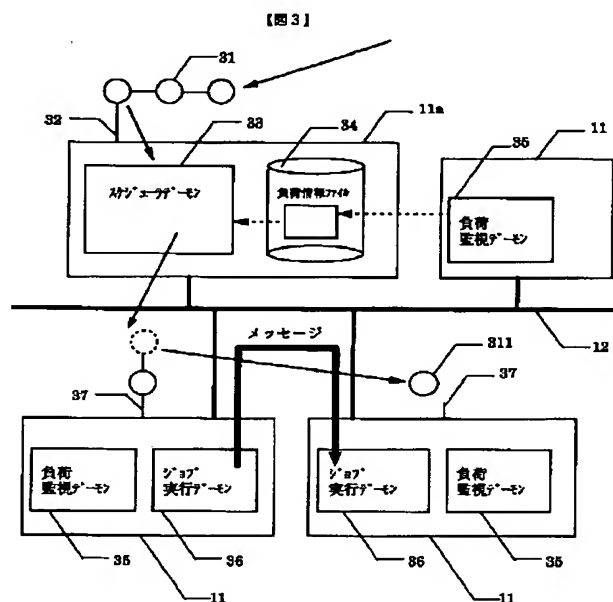
(74) 代理人 弁理士 小川 勝男

(54) 【発明の名称】 負荷分散制御方式及び装置

(57) 【要約】

【課題】 計算機の負荷情報に基づき、スケジューラ計算機がジョブまたはプロセスを計算機に割り当てる場合、時間が経過した負荷情報等が元で誤判断する場合があったが、計算機の正確な負荷情報を元に再割り当てが難しく効果的に計算機間の負荷の不均衡を是正することができなかった。

【解決手段】 任意の計算機間でノード間通信処理を行う際、送信メッセージに自計算機の負荷状態を示す負荷情報を付加するようにすることで、計算機間で負荷情報の交換を図り、負荷の高い計算機に割り当てられた実行待ちジョブまたはプロセスを負荷の低い計算機に移動することによって、効果的に負荷の不均衡を是正する。



**【特許請求の範囲】**

【請求項 1】 ネットワークを介し複数の計算機が接続された並列計算機システムにおいて、ユーザから受け付けたプロセス及びジョブをどの計算機に実行させるべきかを判断するスケジューラデーモンを具備するスケジューラ計算機と、実行依頼されたプロセス及びジョブを単に実行する計算機（以下単に計算機といえちこちを指す）から構成されるシステムにおいて、各計算機間でメッセージパッシング等のノード間通信（計算機間の通信処理）を行う際に、送信メッセージに自計算機の負荷情報を付加することにより計算機間で負荷情報を交換し、それにより計算機間の負荷の均衡を保つことを特長とする負荷分散制御方式。

【請求項 2】 複数の計算機がネットワークで接続されており、複数の計算機間でノード間通信を行う並列システムにおいて、各計算機に自計算機からの送信メッセージに自計算機の負荷情報を付加する装置と、受信メッセージに付加されている他計算機の負荷情報を認知するための装置と、それら負荷情報から負荷の高い計算機に割り当てられた実行待ちジョブまたはプロセスを負荷の低い計算機に移動する装置を具備することによって、効果的に負荷の不均衡を是正することを特長とする負荷分散制御装置。

**【発明の詳細な説明】****【0001】**

【発明の属する技術分野】 本発明は、並列計算機システムにおけるプロセス及びジョブの実行制御に係わり、計算機の負荷分散制御方式に関する。

**【0002】**

【従来の技術】 従来の方法の一例として、特開平 0 9 - 1 6 0 8 8 4 のように共有メモリを介して各計算機が他計算機と負荷情報を交換しあい、計算機間の負荷の均衡を図る方法があげられる。

【0003】 また上記に依らず、計算機の負荷情報を元に、ジョブまたはプロセスのスケジューリングをスケジューラ計算機が行う場合、各計算機からスケジューラ計算機に向けてのみ負荷情報の送信が行われていた。

**【0004】**

【発明が解決しようとする課題】 共有メモリを介して計算機の負荷情報を交換する方法は、共有メモリを有するシステムのみに適用可能であり、それ故適用範囲が限定されていた。

【0005】 また負荷情報を元に、スケジューラ計算機がジョブまたはプロセスをスケジューリングする場合、実行以前に当該ジョブまたはプロセス実行時の計算機に与える負荷を正確に判断することが難しいため、結果的に誤った判断を下してしまう場合があり、それにより計算機間で負荷の不均衡が発生していた。

【0006】 なおこの負荷の不均衡が発生した場合、それを是正するために、負荷の不均衡が発生している計算

機からスケジューラ計算機に負荷情報を送信する方法は、自計算機に更なる負荷を与えるだけでなく、スケジューラ計算機へのネットワークの局所的な負荷やオーバーヘッドの増大を招き、これらの処理に時間を費やすことで送付された負荷情報の正確さが徐々に失われるため、スケジューラ計算機を介した方法では正確な負荷の不均衡是正の実現が困難であった。

**【0007】**

【課題を解決するための手段】 計算機間で負荷の不均衡が発生した場合、それを是正できるようにするため、任意の計算機間でノード間通信処理を行う際、送信メッセージに自計算機の負荷状態を示す負荷情報を付加するようにする。

【0008】 送信メッセージ及び負荷情報を受け取った計算機は、自計算機の負荷と、メッセージ送信側の負荷とを比較し、必要に応じて、負荷（未実行ジョブまたはプロセス）の移動を行うようにする。

【0009】 即ち、メッセージ送信側の計算機の負荷が高く、メッセージ受信側の計算機の負荷が低い場合には、送信メッセージに負荷が高い旨の負荷情報が付加され、メッセージ送信側計算機の未実行ジョブまたはプロセスを、メッセージ受信側計算機に移動するようにする。

【0010】 逆にメッセージ送信側の計算機の負荷が低く、メッセージ受信側の計算機の負荷が高い場合には、送信メッセージに負荷が低い旨の負荷情報が付加され、メッセージ受信側計算機の未実行ジョブまたはプロセスを、メッセージ送信側計算機に移動するようにする。

【0011】 なお負荷の高低は、両計算機間の負荷の相対的な比較に基づくが、負荷の差異が微小である場合には、両計算機がほぼ同一の負荷を有していると見なせるため負荷の移動には及ばない。よってこのように両計算機間の負荷の差異が微小と見なせるか、見なせないかの基準となる、差異許容値  $k$  を導入し、それによりメッセージ受信側計算機は、自計算機の方が負荷は高いか、低いか、または同一かを判断する。

【0012】 また未実行ジョブまたはプロセスを移動する際には、移動の旨及び、両計算機の最新の負荷情報をスケジューラ計算機に伝え、スケジューラ計算機は、当該計算機の負荷情報を入手することにより、それを元にスケジューリングを行うようにする。

**【0013】**

【発明の実施の形態】 以下、本発明の一実施例について図面により詳細に説明する。

【0014】 図 1 は本発明を適用した並列計算機システムのシステム構成を示すブロック図である。11 は計算機を示し、複数台の計算機 11 がネットワーク 12 にそれぞれ接続し、システムを構成する。図 1 では計算機は 4 台を示しているが、計算機 11 の台数は任意である。

【0015】 図 2 は計算機 11 の詳細な構成を示すプロ

は発生しない。

【0035】(b)計算機Sは移動する未実行ジョブまたはプロセス311を自計算機の実行待ちキュー37から切り離し、スケジューラデーモン33から受け取ったジョブまたはプロセスの実行依頼内容を計算機Rに送付し、計算機Rは受け取った実行依頼内容を元に、自計算機の実行待ちキュー37に未実行ジョブまたはプロセス311をキューイングする。

【0036】(d)計算機S及び計算機Rは、共に自計算機の負荷が変化したことを伝えるため、自計算機の負荷をスケジューラデーモン33に報告する。

【0037】(iii)ケースCの場合負荷の移動は発生しない。(図7、図8のケースC参照)

【0038】

【発明の効果】共有メモリを有しない計算機システムにおいても各計算機の正確な負荷情報によりスケジューラ計算機が負荷の低い計算機へジョブまたはプロセスをスケジューリングするためシステム全体のスループットを向上させることができる

【図面の簡単な説明】

【図1】本発明の一実施例を示すシステム構成ブロック図

【図2】図1の構成例に示した計算機の構成図

【図3】本発明の一実施例を示すシステム構成図

【図4】メッセージ送信側計算機とメッセージ受信側計

算機との情報のやりとりを示す図

【図5】メッセージ送信側計算機とメッセージ受信側計算機との情報のやりとりを示す図

【図6】メッセージ送信側計算機とメッセージ受信側計算機との情報のやりとりを示す図

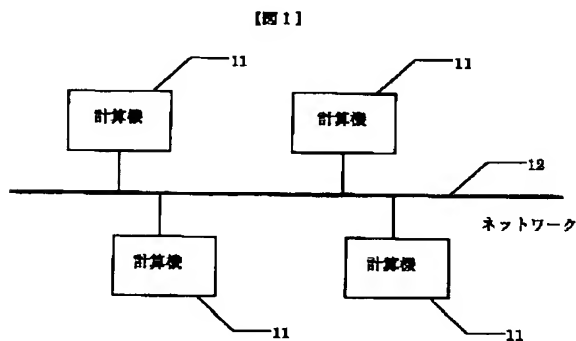
【図7】メッセージ送信側計算機のフローチャート図

【図8】メッセージ受信側計算機のフローチャート図

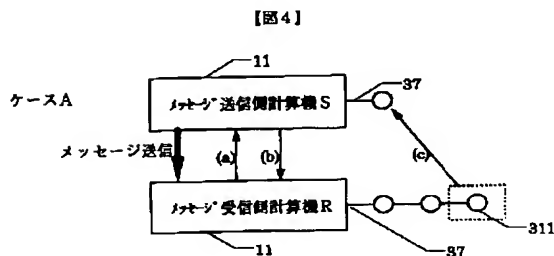
【符号の説明】

- 11：計算機
- 12：ネットワーク
- 21：CPU
- 22：メモリ
- 23：外部記憶装置
- 231：負荷情報ファイル
- 24：バス
- 31：ユーザから実行依頼されたジョブまたはプロセス
- 311：ユーザから実行依頼されたジョブまたはプロセスで、スケジューラ計算機から計算機を割り当てられたもの
- 32：ジョブまたはプロセス受け付けキュー
- 33：スケジューラデーモンプログラム
- 34：負荷情報ファイル
- 35：負荷監視デーモン
- 36：ジョブ実行デーモン
- 37：実行待ちプロセスまたはジョブキュー

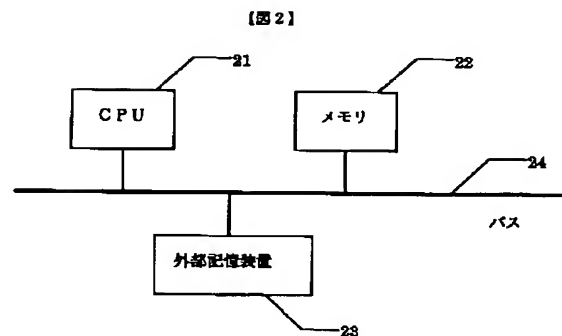
【図1】



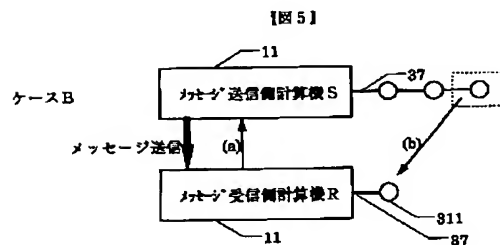
【図4】



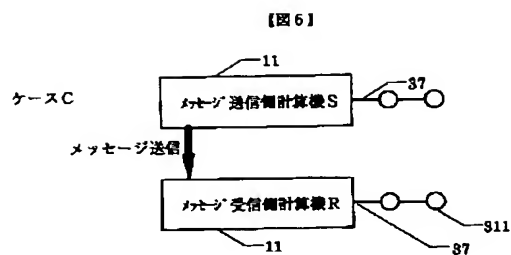
【図2】



【図5】



【图 6】



【图 8】

